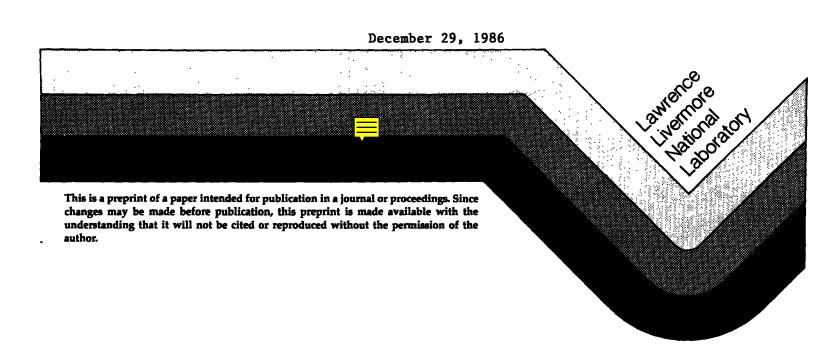
UCRL-95905 PREPRINT

SIRCULATION COPY SUBJECT TO RECALL TWO WEEKS

The Numerical Solution of Higher Index
Differential/Algebraic Equations by Implicit
Runge-Kutta Methods

Kathryn E. Brenan UCLA

Linda R. Petzold Lawrence Livermore National Laboratories



DISCLAIMER

This document was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial products, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government thereof, and shall not be used for advertising or product endorsement purposes.

The Numerical Solution of Higher Index Differential/Algebraic Equations by Implicit Runge-Kutta Methods

Kathryn E. Brenan* UCLA

Linda R. Petzold[†]
Lawrence Livermore National Laboratories

January 5, 1987

Abstract

In this paper we study the order, stability, and convergence properties of implicit Runge-Kutta (IRK) methods applied to differential/algebraic systems with index greater than one. These methods do not in general attain the same order of accuracy for higher index differential/algebraic systems as they do for index one systems or for purely differential systems. We derive necessary and sufficient conditions on the method coefficients to ensure that the local and global errors of the method attain a given order of accuracy for high index linear constant coefficient systems. We study IRK methods applied to nonlinear semi-explicit index two systems and derive a sufficient set of conditions which ensure that a method is accurate to a given order for these systems. Finally, we present some numerical experiments which illustrate these results.

^{*}This work was supported in part by NSF Grant DMS-85-03294.

[†]This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract No. W-7405-Eng-48.

1 Introduction

In this paper we extend the results for order, stability, and convergence of implicit Runge-Kutta (IRK) methods derived for index one systems by Petsold [17] to higher index differential/algebraic systems. It is well known that these methods often do not attain the same order of accuracy for differential/algebraic systems as they do for purely differential systems. Petzold has studied their behavior on uniformly index one systems of the form,

$$F(y,y',t)=0 (1.1)$$

when consistent initial values $y(t_0)$ are given. We examine two classes of DAEs not considered in [17], to understand how the order of accuracy of an IRK method depends on the index of the DAE system as well as the method coefficients. First we study solvable linear constant coefficient systems

$$Ay' + By = g(t), \tag{1.2}$$

of arbitrary index ν , where A and B are square constant matrices and g(t) is a smooth function. Then we develop a convergence theory for IRK methods applied to nonlinear semi-explicit index two systems,

$$f(x, x', y, t) = 0$$

 $g(x, y, t) = 0,$ (1.3)

where $(\partial f/\partial x')^{-1}$ exists and is bounded in some neighborhood of the solution and $\partial g/\partial y$ has constant rank.

We formally apply an M-stage IRK method to a DAE (1.1) to obtain the system of difference equations

$$F(y_{n-1} + h \sum_{j=1}^{M} a_{ij} Y_j', Y_i', t_{n-1} + c_i h) = 0 \quad i = 1, 2, ..., M$$

$$y_n = y_{n-1} + h \sum_{i=1}^{M} b_i Y_i'$$
(1.4)

where $h = t_n - t_{n-1}$. We will assume throughout this paper that the coefficient matrix $A = (a_{ij})$ of the Runge-Kutta method is nonsingular. Note that this method reduces to a standard IRK method when applied to a system of explicit ordinary differential equations (ODEs).

In section 2 we study linear constant coefficient systems of arbitrary index ν . We derive necessary and sufficient conditions on the method coefficients to ensure that the local error of the method attains a given order of accuracy for these systems. We also investigate the error propagation properties of IRK methods applied to these systems, and derive an expression for the global error.

In section 3 we develop a convergence theory for IRK methods applied to nonlinear semi-explicit index two systems of the form (1.3). We derive a sufficient set of conditions which ensure that a method is accurate to a given order for these systems.

In the last section we describe some numerical experiments which illustrate the order reduction effects predicted by the theory, and also raise some interesting questions for future research.

2 Linear Constant Coefficient Systems

In this section we derive conditions that are necessary and sufficient to ensure that the local error of an implicit Runge-Kutta method attains a given order when applied to linear constant coefficient systems of arbitrary index ν . Then we study error propagation for constant coefficient higher index systems, and derive an expression for the global error.

Consider the linear constant coefficient DAE (1.2)

$$Ay' + By = g(t) \tag{2.1}$$

of index ν . We assume this system is solvable, so that there exist nonsingular matrices P and Q which decouple the system [10],

$$PAQ = \begin{pmatrix} I & 0 \\ 0 & N \end{pmatrix} \qquad PBQ = \begin{pmatrix} C & 0 \\ 0 & I \end{pmatrix}$$
 (2.2)

where I is an identity matrix and N is a block diagonal matrix, $N = \operatorname{diag}(N_1, N_2, \ldots, N_L)$ composed of blocks of the form

$$N_i = \begin{pmatrix} 0 & & & \\ 1 & 0 & & \\ & \ddots & \ddots & \\ & & 1 & 0 \end{pmatrix}. \tag{2.3}$$

Applying the IRK method to (2.1), we have

$$AY'_{i} + B(y_{n-1} + h \sum_{j=1}^{M} a_{ij}Y'_{j}) = g(t_{n-1} + c_{i}h) \quad i = 1, 2, ..., M$$

$$y_{n} = y_{n-1} + h \sum_{i=1}^{M} b_{i}Y'_{i}. \tag{2.4}$$

Premultiplying these difference equations by P and letting $\tilde{y}_n = Q^{-1}y_n$, $\tilde{Y}'_i = Q^{-1}Y'_i$, $\tilde{g}(t) = Pg(t)$, we obtain

$$(PAQ)\tilde{Y}'_{i} + (PBQ)(\tilde{y}_{n-1} + h\sum_{j=1}^{M} a_{ij}\tilde{Y}'_{j}) = \tilde{g}(t_{n-1} + c_{i}h) \quad i = 1, 2, ..., M$$

$$\tilde{y}_{n} = \tilde{y}_{n-1} + h\sum_{i=1}^{M} b_{i}\tilde{Y}'_{i}.$$

Note that the differential and algebraic parts of the system are decoupled in this form. In addition, the algebraic subsystems are decoupled from one another. Thus it is sufficient to study the behavior of the IRK method on a canonical algebraic subsystem to understand its behavior on general linear constant coefficient systems.

Consider then a canonical algebraic subsystem of index ν

$$Ny' + y = g(t) \tag{2.5}$$

where N is a $\nu \times \nu$ matrix of the form (2.3), $g(t) = (g_1(t), g_2(t), \dots, g_{\nu}(t))^T$, and $g(t) = (y_1(t), y_2(t), \dots, y_{\nu}(t))^T$. The solution to (2.5) is given by

$$y_1(t) = g_1(t)$$

$$y_2(t) = g_2(t) - g_1'(t)$$

$$\vdots$$

$$y_{\nu}(t) = g_{\nu}(t) + \sum_{i=1}^{\nu-1} (-1)^{\nu-i} g_i^{(\nu-i)}(t).$$

Applying the IRK method to (2.5), we obtain

$$NY_{i}' + \left(y_{n-1} + h \sum_{j=1}^{M} a_{ij} Y_{j}'\right) = g(t_{n-1} + c_{i}h) \qquad i = 1, 2, ..., M$$

$$y_{n} = y_{n-1} + h \sum_{i=1}^{M} b_{i} Y_{i}'. \qquad (2.6)$$

Let $Y_i' = (Y_{i,1}', Y_{i,2}', \ldots, Y_{i,\nu}')^T$ where $Y_{i,j}'$ denotes the *i*th stage derivative corresponding to the *j*th component of the solution vector, namely the 'index *j*' variable $y_j(t)$. Because the coefficient matrix $A = (a_{ij})$ of the IRK method is nonsingular, the difference equations (2.6) can be solved uniquely for the stage derivatives $Y' = (Y_1', Y_2', \ldots, Y_M')^T$. Because of the structure of N, the solution to the *i*th equation in the original system (2.5) depends only on the solutions to the first (i-1) equations. A similar dependency is present in the difference equations (2.6), allowing us to solve first for the stage derivatives $Y_{1,1}', Y_{2,1}', \ldots, Y_{M,1}'$ corresponding to the index one variable $y_1(t)$, second for the stage derivatives $Y_{1,2}', Y_{2,2}', \ldots, Y_{M,2}'$ corresponding to the index two variable $y_2(t)$, etc. This analysis has already been done in [17] for the index one variable. Here we extend the analysis to more general index ν systems.

In general, for each component of the solution $y_j(t)$, $j = 1, 2, ..., \nu$, we solve a subset of M equations from system (2.6) for the corresponding stage derivatives:

$$(Y'_{1,j}, Y'_{2,j}, \dots, Y'_{M,j})^{T} = (1/h) A^{-1} ((g_{j}(t_{n-1} + c_{1}h), g_{j}(t_{n-1} + c_{2}h), \dots, g_{j}(t_{n-1} + c_{M}h))^{T} -\epsilon_{M} y_{n-1,j} - (Y'_{1,j-1}, Y'_{2,j-1}, \dots, Y'_{M,j-1})^{T})$$

$$(2.7)$$

where $\epsilon_M = (1, 1, ..., 1)^T$. Note that the stage derivatives depend only on the numerical solution $y_{n-1,j}$ and on the stage derivatives of the (j-1)st variable. Define

$$G_{i} = \begin{pmatrix} (g_{i}(t_{n-1} + c_{1}h) - g_{i}(t_{n-1}))/h \\ (g_{i}(t_{n-1} + c_{2}h) - g_{i}(t_{n-1}))/h \\ \vdots \\ (g_{i}(t_{n-1} + c_{M}h) - g_{i}(t_{n-1}))/h \end{pmatrix}$$

$$c^{i} = (c_{1}^{i}, c_{2}^{i}, \dots, c_{M}^{i})^{T}.$$

Utilizing the local error assumption,

$$y_{n-1} = y(t_{n-1}) (2.8)$$

and substituting for the (j-1)st stage derivatives in (2.7), we find the following expressions for the *local* stage derivatives:

$$(Y'_{1,1}, Y'_{2,1}, \dots, Y'_{M,1})^{T} = A^{-1}G_{1}$$

$$(Y'_{1,2}, Y'_{2,2}, \dots, Y'_{M,2})^{T} = A^{-1}G_{2} - (1/h)A^{-2}G_{1} + (1/h)A^{-1}\epsilon_{M}g'_{1}(t_{n-1})$$

$$(Y'_{1,3}, Y'_{2,3}, \dots, Y'_{M,3})^{T} = A^{-1}G_{3} + (1/h)A^{-1}\epsilon_{M}[g'_{2}(t_{n-1}) - g''_{1}(t_{n-1})]$$

$$- (1/h)A^{-2}G_{2} + (1/h^{2})A^{-3}G_{1}$$

$$- (1/h^{2})A^{-2}\epsilon_{M}g'_{1}(t_{n-1})$$

$$(2.9)$$

and similar expressions for the remaining higher index variables. We define the local error d_n by

$$d_n = y(t_{n-1}) + h \sum_{i=1}^{M} b_i Y_i' - y(t_n)$$
 (2.10)

where $d_n = (d_{n,1}, d_{n,2}, \ldots, d_{n,\nu})^T$ and Y_i' represent the local stage derivatives given by (2.9). Expanding (2.10) in a Taylor series about t_{n-1} , and equating like powers of h, it is easy to see as in [17] that the local error $d_{n,1}$ in the index one variable satisfies

$$d_{n,1} = O(h^{k_{a,1}+1}), (2.11)$$

when $b^T A^{-1} c^i = 1$ for $i = 1, 2, ..., k_{a,1}$. $k_{a,1}$ is the algebraic order of the IRK method applied to index one constant coefficient systems. For the index two variable, the local error $d_{n,2}$ is given by

$$d_{n,2} = y_2(t_{n-1}) - y_2(t_n) + hb^T A^{-1} G_2 + b^T A^{-1} \epsilon_M g_1'(t_{n-1}) - b^T A^{-2} G_1$$
 (2.12)

Note that $g_2(t_{n-1}) - g_2(t_n) + hb^T A^{-1}G_2 = O(h^{k_{a,1}+1})$ if we assume the IRK method has algebraic order $k_{a,1}$ on index one problems. Then expand the remaining terms in (2.12) in a Taylor series about t_{n-1} and equate like powers of h to obtain the following set of order conditions for index two constant coefficient systems,

$$b^T A^{-1} \epsilon_M = b^T A^{-2} c^1$$

 $b^T A^{-2} c^i = i, i = 2, 3, ..., k_{a,2}.$

We define the algebraic order of the IRK method applied to index two constant coefficient canonical systems to be $k_{a,2}$ if these conditions are satisfied. The local error for a general index two constant coefficient system thus satisfies $d_{a,2} = O(h^{k_{a,2}}) + O(h^{k_{a,1}+1})$.

This analysis can be extended in a straightforward manner to the general index ν case. For completeness we list here the additional algebraic order conditions and the corresponding asymptotic behavior of the local error for both the index three case and the most general case. For an index three system,

$$d_{n,3} = O(h^{k_{a,3}-1}) + O(h^{k_{a,2}}) + O(h^{k_{a,1}+1})$$

where $k_{a,3}$ is the largest integer such that

$$b^{T} A^{-2} \epsilon_{M} = b^{T} A^{-3} c^{1}$$

$$b^{T} A^{-1} \epsilon_{M} = b^{T} A^{-3} c^{2} / 2$$

$$b^{T} A^{-3} c^{i} = i(i-1), \quad i = 3, 4, \dots, k_{a,3}.$$

Finally, for a general constant coefficient index ν system, the local error satisfies

$$d_{n,\nu} = O(h^{k_{a,\nu}-\nu+2}) + O(h^{k_{a,\nu-1}-\nu+3}) + \dots + O(h^{k_{a,1}+1})$$
 (2.13)

where $k_{a,\nu}$ is the largest integer such that

$$b^{T} A^{-i} \epsilon_{M} = b^{T} A^{-\nu} c^{\nu-i} / (\nu - i)! \quad i = 1, 2, ..., \nu - 1$$

$$b^{T} A^{-\nu} c^{i} = i(i-1) ... (i-\nu+1), \quad i = \nu, \nu+1, ..., k_{a,\nu}.$$

Clearly, the higher the index the more difficult it is to find IRK methods which are convergent in all the variables. Finding an IRK method having the same rate of convergence in all of the variables similarly poses severe restrictions on the coefficients.

Next we examine the propagation of errors for IRK methods applied to linear constant coefficient systems. Consider solving (2.5) by the perturbed Runge-Kutta method,

$$NZ'_{i} + (z_{n-1} + h \sum_{j=1}^{M} a_{ij}Z'_{j} - \delta_{n}^{(i)}) = g(t_{n-1} + c_{i}h) \quad i = 1, 2, ..., M$$

$$z_{n} = z_{n-1} + h \sum_{i=1}^{M} b_{i}Z'_{i} - \delta_{n}^{(M+1)}, \qquad (2.14)$$

where the perturbations $\delta_n^{(i)} = (\delta_{n,1}^{(i)}, \delta_{n,2}^{(i)}, \dots, \delta_{n,\nu}^{(i)})^T$ satisfy $\|\delta_n^{(i)}\| \leq \Delta$ for $i = 1, 2, \dots, M + 1$. The perturbations could be due to roundoff error, errors in solving the linear systems at each stage, or could be interpreted as truncation errors at each stage (see section 3). Subtracting (2.14) from (2.6), and defining $e_n = y_n - z_n$, $E_i' = Y_i' - Z_i'$, we obtain an expression for the difference between these two solutions

$$NE'_{i} + \left(e_{n-1} + h \sum_{j=1}^{M} a_{ij} E'_{j} + \delta_{n}^{(i)}\right) = 0, \quad i = 1, 2, ..., M$$

$$e_{n} = e_{n-1} + h \sum_{i=1}^{M} b_{i} E'_{i} + \delta_{n}^{(M+1)}. \tag{2.15}$$

By solving the first equation in (2.15) for E'_i and substituting into the second equation, we can obtain a relation describing the error propagation

of the method. For linear constant coefficient index one systems this was done in Petsold [17] and resulted in

$$e_{n,1} = (1 - b^T A^{-1} \epsilon_M) e_{n-1,1} - (b^T A^{-1} \delta_{n,1} - \delta_{n,1}^{(M+1)}), \qquad (2.16)$$

where $\delta_{n,j} = (\delta_{n,j}^{(1)}, \delta_{n,j}^{(2)}, \dots, \delta_{n,j}^{(M)})^T$. The recurrence (2.16) is unstable unless $|1 - b^T A^{-1} \epsilon_M| < 1$. Hence we will require as in [17] that the IRK method satisfy the strict stability condition

$$\mid 1 - b^T A^{-1} \epsilon_M \mid < 1. \tag{2.17}$$

The error propagation relation for the index two variable is given by

$$e_{n,2} = (1 - b^T A^{-1} \epsilon_M) e_{n-1,2} - (b^T A^{-1} \delta_{n,2} - \delta_{n,2}^{(M+1)}) + (1/h) b^T A^{-2} (\delta_{n,1} + \epsilon_M e_{n-1,1}),$$
 (2.18)

while for the index three variable it is

$$e_{n,3} = (1 - b^T A^{-1} \epsilon_M) e_{n-1,3} - (b^T A^{-1} \delta_{n,3} - \delta_{n,3}^{(M+1)}) + (1/h) b^T A^{-2} (\delta_{n,2} + \epsilon_M e_{n-1,2}) - (1/h^2) b^T A^{-3} (\delta_{n,1} + \epsilon_M e_{n-1,1}).$$
(2.19)

Finally, for the general index ν case, the stability relation can be shown to be

$$e_{n,\nu} = (1 - b^{T} A^{-1} \epsilon_{M}) e_{n-1,\nu} - (b^{T} A^{-1} \delta_{n,\nu} - \delta_{n,\nu}^{(M+1)}) - \sum_{i=1}^{\nu-1} \frac{(-1)^{i}}{h^{i}} b^{T} A^{-i-1} (\delta_{n,\nu-i} + \epsilon_{M} \epsilon_{n-1,\nu-i}).$$
 (2.20)

Note that the strict stability condition is no longer sufficient to insure stability, in a strict mathematical sense, of the IRK method when applied to linear constant coefficient systems of index greater than one. For small stepsizes, roundoff errors can be significant for higher index variables.

These methods can be useful for the solution of higher index systems, provided that we understand the implications of the error propagation relations given above. We can see that the sensitivity, to roundoff errors is confined to the higher index variables of the system, and does not propagate back into the lower index variables. This observation holds also for the nonlinear semi-explicit index two systems that we study in the next section. Finally, using the error propagation relations above, we can extend the conclusions of Petzold [17] for global error in solving linear constant coefficient index one systems to higher index systems as follows.

Definition 2.1 The constant coefficient order of an IRK method (1.4) is equal to $k_{c,\nu}$ if the method converges with global error $O(h^{k_{c,\nu}})$ for all solvable linear constant coefficient systems (1.2) of index $\leq \nu$.

Theorem 2.1 Suppose the IRK method (1.4) satisfies the strict stability condition. Then the constant coefficient order $k_{c,\nu}$ of the global error of this method is given by

$$k_{c,\nu} = \min_{1 \le i \le \nu} (k_d, k_{a,i} - \nu + 2)$$
 (2.21)

where k_d is the order of the method for purely differential (nonstiff) systems.

Finally, we present some results on the order of accuracy of some IRK methods from the stiff ODE literature applied to index one and index two linear constant coefficient systems. We have chosen to investigate these particular methods because our numerical experience [2] with IRK methods applied to DAEs has led us to conclude that it is very desirable for a method to be L-stable, or even better to be stiffly accurate, and also because these methods can be implemented efficiently. One reason why L-stable methods appear promising is that they perform very well when applied to index one and semi-explicit index two and index three systems, even when the initial values contain small errors. Recall that a method is L-stable if it is Astable and if $\lim_{R_0(h\lambda)\to-\infty} |y_{n+1}/y_n|=0$, when applied to the test problem $y' = \lambda y$. For IRK methods, this condition is equivalent to requiring that $|1 - b^T A^{-1} \epsilon_M| = 0$. Stiffly accurate methods [18]) are L-stable methods which satisfy the additional requirement that $c_M = 1$, $a_{Mj} = b_j$ for j =1, 2, ..., M. Thus $b^T A^{-1} = (0, 0, ..., 0, 1)^T$ for stiffly accurate methods. The L-stable methods we have chosen to investigate here are:

- (1) 2-stage, '2nd order' Singly Implicit method (SIRK) [4], with $\lambda = 1 \sqrt{2}/2$
- (2) 5-stage, '4th order' Diagonally Implicit method (DIRK) [1] [7]
- (3) 3-stage, '3rd order' Singly Implicit method (SIRK) [4], with $1/\lambda$ the root of the Laguerre polynomial of degree three
- (4) 7-stage, '3rd order' Extrapolation method based on fully implicit backward Euler and polynomial extrapolation, written as a semi-implicit Runge-Kutta method

Methods (1) and (2) are stiffly accurate. The results are given in Table 1, where it can be seen that, as observed above, it is difficult to maintain the same rate of convergence in all of the variables for linear constant coefficient index two systems.

Table 2.1: Order of Consistency for LCC Canonical Systems

L-Stable Methods	ODE order kd	Index 1 Order ka,1	Index 2 Order ka,2
1. Two-stage SIRK	$O(h^2)$	Exact	$O(h^2)$
2. Five-stage DIRK	$O(h^4)$	Exact	O(h)
3. Three-stage SIRK	$O(h^3)$	$O(h^3)$	$O(h^2)$
4. Seven-stage Extrp.	$O(h^3)$	Exact	$O(h^3)$

3 Semi-Explicit Nonlinear Index Two Systems

In this section we study nonlinear semi-explicit index two systems of the form

$$f(x, x', y, t) = 0 g(x, y, t) = 0, (3.1)$$

where we will assume that $(\partial f/\partial x')^{-1}$ exists and is bounded in some neighborhood of the solution, $\partial g/\partial y$ has constant rank, and f and g have as many continuous partial derivatives as desired in a neighborhood of the solution. We give a set of order conditions which are sufficient to ensure that a method is accurate to a given order for these systems.

To state our results, we first recall the definitions of internal order and internal local truncation error given in [17].

Definition 3.1 The ith internal local truncation error $\delta_i^{(n)}$ at t_n of an M-stage implicit Runge-Kutta method (1.4) is given by

$$\delta_i^{(n)} = x(t_{n-1}) + h \sum_{j=1}^M a_{ij} x'(t_{n-1} + c_j h) - x(t_{n-1} + c_i h), \quad i = 1, \ldots, M,$$

$$\delta_{M+1}^{(n)} = x(t_{n-1}) + h \sum_{i=1}^{M} b_i x'(t_{n-1} + c_i h) - x(t_n). \tag{3.2}$$

Definition 3.2 The internal order k_I of an M-stage implicit Runge-Kutta method (1.4) is given by

$$k_I = \min(k_1, \ldots, k_M, k_{M+1})$$

where

$$\delta_i^{(n)} = O(h^{k_i+1}), \quad i = 1, \ldots, (M+1).$$

As in [17] [9], it is simple to find the internal order of an implicit Runge-Kutta method in terms of its coefficients by expanding (3.2) in Taylor series around t_{n-1} , leading to the result that the internal order of an M-stage implicit Runge-Kutta method is equal to k_I iff the method coefficients satisfy

$$\sum_{j=1}^{M} a_{ij} c_j^{k-1} = \frac{c_i^k}{k}, \quad i = 1, \dots, M,$$

$$\sum_{j=1}^{M} b_j c_j^{k-1} = \frac{1}{k}$$

for $k=1,\ldots,k_I$.

Then we can state the following result.

Theorem 3.1 Given the nonlinear, semi-explicit index two system (3.1) to be solved numerically by the M-stage IRK method (1.4), suppose

- 1) The IRK method has internal order k1
- 2) The IRK method satisfies the strict stability condition
- 3) The initial conditions satisfy $||(I-H)e_0^x|| = O(h^{k_G})$, where I-H is a projection operator defined below, $e_0^x = x_0 x(t_0)$ and $k_G = \min(k_d, k_I + 1)$

Then the global errors in the numerical solution x_n and y_n are $O(h^{k_0})$ and $O(h^{k_1})$, respectively.

Proof. We will first prove this theorem for the simpler index 2 system

$$x' + g_1(x, y, t) = 0$$

 $g_2(x, t) = 0,$ (3.3)

where $[(\partial g_2/\partial x)(\partial g_1/\partial y)]^{-1}$ exists and is bounded in a neighborhood of the solution, and then show that the results extend to the more general system (3.1).

Consider the M-stage IRK applied to (3.3):

$$X'_{i} + g_{1} \left(x_{n-1} + h \sum_{j=1}^{M} a_{ij} X'_{j}, y_{n-1} + h \sum_{j=1}^{M} a_{ij} Y'_{j}, t_{i} \right) = 0$$

$$g_{2} \left(x_{n-1} + h \sum_{j=1}^{M} a_{ij} X'_{j}, t_{i} \right) = 0 \qquad i = 1, 2, ..., M$$

$$x_{n} = x_{n-1} + h \sum_{i=1}^{M} b_{i} X'_{i} \qquad ,$$

$$y_{n} = y_{n-1} + h \sum_{j=1}^{M} b_{j} Y'_{j}, \qquad (3.4)$$

where $t_i = t_{n-1} + c_i h$. It is convenient to define intermediate stage values for x and y at t_i :

$$X_i = x_{n-1} + h \sum_{j=1}^M a_{ij} X_j'$$

$$Y_i = y_{n-1} + h \sum_{j=1}^{M} a_{ij} Y_j'.$$

The true solution satisfies

$$x'(t_i) + g_1(x(t_i), y(t_i), t_i) = 0$$

 $g_2(x(t_i), t_i) = 0$ $i = 1, 2, ..., M$

$$x(t_n) = x(t_{n-1}) + h \sum_{i=1}^{M} b_i x'(t_{n-1} + c_i h) - \delta_{M+1}^{x(n)}$$

$$y(t_n) = y(t_{n-1}) + h \sum_{i=1}^{M} b_i y'(t_{n-1} + c_i h) - \delta_{M+1}^{y(n)}, \qquad (3.5)$$

where

$$egin{split} x(t_i) &= x(t_{n-1}) + h \sum_{j=1}^M a_{ij} x'(t_{n-1} + c_j h) - \delta_i^{x(n)} \ y(t_i) &= y(t_{n-1}) + h \sum_{j=1}^M a_{ij} y'(t_{n-1} + c_j h) - \delta_i^{y(n)} \quad i = 1, 2, \dots, M. \end{split}$$

Let $G_{21}(t_i) = \partial g_2/\partial x$, $G_{12}(t_i) = \partial g_1/\partial y$, and $G_{11}(t_i) = \partial g_1/\partial x$, where the partial derivatives are evaluated along the true solution at t_i . Subtracting (3.5) from (3.4), we obtain

$$E_i^{xi} + G_{11}(t_i)E_i^x + G_{12}(t_i)E_i^y = \eta_i^x$$

 $G_{21}(t_i)E_i^x = \eta_i^y$ $i = 1, 2, ..., M$

$$e_n^x = e_{n-1}^x + h \sum_{i=1}^M b_i E_i^{xi} + \delta_{M+1}^{x(n)}$$

$$e_n^y = e_{n-1}^y + h \sum_{i=1}^M b_i E_i^{yi} + \delta_{M+1}^{y(n)},$$
(3.6)

where

$$E_{i}^{y} = e_{n-1}^{y} + h \sum_{j=1}^{M} a_{ij} E_{j}^{yi} + \delta_{i}^{y(n)}$$

$$E_{i}^{z} = e_{n-1}^{z} + h \sum_{i=1}^{M} a_{ij} E_{j}^{zi} + \delta_{i}^{z(n)}, \qquad (3.7)$$

and $E_i^{xi} = X_i' - x'(t_i)$, $E_i^{yi} = Y_i' - y'(t_i)$, $E_i^y = Y_i - y(t_i)$, $E_i^x = X_i - x(t_i)$, $e_n^x = x_n - x(t_n)$, and $e_n^y = y_n - y(t_n)$. The η_i terms are the sum of residuals from the Newton iteration and higher order terms in E_i^z and E_i^y .

We can eliminate E_i^y in terms of $E_i^{x'}$ by multiplying the first equation in (3.6) by $G_{21}(t_i)$ and solving for E_i^y ,

$$E_i^y = -M_i E_i^{x'} - M_i G_{11}(t_i) E_i^x + M_i \eta_i^x,$$
 (3.8)

where $M_i = (G_{21}(t_i)G_{12}(t_i))^{-1}G_{21}(t_i)$.

Let $H_i = G_{12}(t_i)M_i$. Multiply the first equation in (3.6) by $I - H_i$ and substitute (3.8) for E_i^y . Multiply the second equation in (3.6) by $G_{12}(t_i)(G_{21}(t_i)G_{12}(t_i))^{-1}$. Equations (3.6) may now be written as

$$(I - H_i)E_i^{x'} + N_iE_i^x = \tilde{\eta}_i^x$$

$$H_iE_i^x = \tilde{\eta}_i^y \qquad i = 1, 2, \dots, M$$

$$e_{n}^{x} = e_{n-1}^{x} + h \sum_{i=1}^{M} b_{i} E_{i}^{x'} + \delta_{M+1}^{x(n)}$$

$$e_{n}^{y} = e_{n-1}^{y} + h \sum_{i=1}^{M} b_{i} E_{i}^{y'} + \delta_{M+1}^{y(n)},$$
(3.9)

where $\tilde{\eta}_i^y = G_{12}(t_i)(G_{21}(t_i)G_{12}(t_i))^{-1}\eta_i^y$, $\tilde{\eta}_i^z = (I - H_i)\eta_i^z$ and $N_i = (I - H_i)G_{11}(t_i).$

Define

$$\begin{split} \tilde{E}_{i}^{x'} &= (I - H_{i}) E_{i}^{x'} & \tilde{E}_{i}^{x'} = H_{i} E_{i}^{x'} \\ \tilde{\delta}_{i}^{z(n)} &= (I - H_{i}) \delta_{i}^{z(n)} & \tilde{\delta}_{i}^{z(n)} = H_{i} \delta_{i}^{x(n)} \\ \tilde{e}_{n}^{z} &= (I - H_{n}) e_{n}^{z} & \tilde{e}_{n}^{z} = H_{n} e_{n}^{z}. \end{split}$$

Rewrite the first equation in (3.9),

$$(I-H_i)E_i^{x'}+N_i\bigg((I-H_i)E_i^x+H_iE_i^x\bigg)=\tilde{\eta}_i^x.$$

Substituting the definition of E_i^x given in equation (3.7) and using the second equation in (3.9), we have an expression for $\tilde{E}_{i}^{z'}$,

$$\tilde{E}_{i}^{x'} + N_{i} \left((I - H_{i}) e_{n-1}^{x} + h \sum_{j=1}^{M} a_{ij} \tilde{E}_{j}^{x'} + h \sum_{j=1}^{M} a_{ij} (H_{j} - H_{i}) E_{j}^{x'} + \tilde{\delta}_{i}^{x(n)} \right) \\
= \tilde{\eta}_{i}^{x} - N_{i} \tilde{\eta}_{i}^{y}.$$
(3.10)

To find an expression for $\tilde{\tilde{E}}_i^{x'}$, from the second equation in (3.9) we have

$$H_i e_{n-1}^x + h \sum_{j=1}^M a_{ij} H_i E_j^{x'} + H_i \delta_i^{x(n)} = \tilde{\eta}_i^y.$$

Thus,

$$H_{i}e_{n-1}^{x} + h\sum_{j=1}^{M} a_{ij}\tilde{\tilde{E}}_{j}^{x'} + h\sum_{j=1}^{M} a_{ij}(H_{i} - H_{j})E_{j}^{x'} + \tilde{\tilde{\delta}}_{i}^{x(n)} = \tilde{\eta}_{i}^{y}.$$
 (3.11)

Now we can rewrite (3.10) and (3.11) noting that $E_j^{x'} = \tilde{E}_j^{x'} + \tilde{\tilde{E}}_j^{x'}$ and $\delta_i^{x(n)} = \tilde{\delta}_i^{x(n)} + \tilde{\tilde{\delta}}_j^{x(n)}$, to obtain

$$\tilde{E}_{i}^{x'} + h \sum_{j=1}^{M} a_{ij} N_{i} \tilde{E}_{j}^{x'} + N_{i} (I - H_{i}) \left(e_{n-1}^{x} + \delta_{i}^{x(n)} \right)
+ h \sum_{j=1}^{M} a_{ij} N_{i} (H_{j} - H_{i}) (\tilde{E}_{j}^{x'} + \tilde{\tilde{E}}_{j}^{x'}) = \tilde{\eta}_{i}^{x} - N_{i} \tilde{\eta}_{i}^{y}
h \sum_{j=1}^{M} a_{ij} \tilde{\tilde{E}}_{j}^{x'} + H_{i} \left(e_{n-1}^{x} + \delta_{i}^{x(n)} \right)
+ h \sum_{i=1}^{M} a_{ij} (H_{i} - H_{j}) (\tilde{E}_{j}^{x'} + \tilde{\tilde{E}}_{j}^{x'}) = \tilde{\eta}_{i}^{y},$$
(3.12)

for i = 1, 2, ..., M. Using their Taylor series expansions about the true solution at t_{n-1} , express N_i and H_i as $N_i = N + O(h)$ and $H_i = H + O(h)$, where N and H are evaluated along the true solution at time t_{n-1} . Then rewrite (3.12) in matrix notation, to obtain:

$$\begin{pmatrix} T_1 & h^2 T_2 \\ h^2 T_3 & h T_4 \end{pmatrix} \begin{pmatrix} \tilde{E}^{x'} \\ \tilde{E}^{x'} \end{pmatrix} = -\begin{pmatrix} S_1 & 0 \\ 0 & S_4 \end{pmatrix} \begin{pmatrix} \mathbf{e}^x_{n-1} + \delta^{x(n)} \\ \mathbf{e}^x_{n-1} + \delta^{x(n)} \end{pmatrix} + \begin{pmatrix} \bar{\eta}^x \\ \bar{\eta}^y \end{pmatrix} \quad (3.13)$$

where

$$\begin{split} \tilde{E}^{z'} &= (\tilde{E}_1^{z'}, \tilde{E}_2^{z'}, \dots, \tilde{E}_M^{z'})^T \\ \tilde{\tilde{E}}^{z'} &= (\tilde{\tilde{E}}_1^{z'}, \tilde{\tilde{E}}_2^{z'}, \dots, \tilde{\tilde{E}}_M^{z'})^T \\ \bar{\tilde{\eta}}^z &= (\tilde{\eta}_1^z - N_1 \tilde{\eta}_1^y, \tilde{\eta}_2^z - N_2 \tilde{\eta}_2^y, \dots, \tilde{\eta}_M^z - N_M \tilde{\eta}_M^y)^T \end{split}$$

$$\begin{split} \tilde{\eta}^{y} &= (\tilde{\eta}_{1}^{y}, \tilde{\eta}_{2}^{y}, \dots, \tilde{\eta}_{M}^{y})^{T} \\ \mathbf{e}_{n-1}^{x} &= (e_{n-1}^{x}, e_{n-1}^{x}, \dots, e_{n-1}^{x})^{T} \\ \delta^{x(n)} &= (\delta_{1}^{x(n)}, \delta_{2}^{x(n)}, \dots, \delta_{M}^{x(n)})^{T}. \end{split}$$

The matrices in (3.13) are given by

$$T_1 = \hat{T}_1 + O(h^2)$$

$$T_4 = \hat{T}_4 + O(h)$$

$$S_1 = \hat{S}_1 + O(h)$$

$$S_4 = \hat{S}_4 + O(h),$$

where $\hat{T}_1 = I_{Md} + hA \otimes N$, $\hat{T}_4 = A \otimes I_d$, $\hat{S}_1 = I_M \otimes (N(I-H))$, $\hat{S}_4 = I_M \otimes H$, T_2 and T_3 are O(1), and d is the dimension of x in (3.3). Here for clarity, we have denoted the dimensions of the identity matrices by subscripts. However, since the matrix I_d occurs frequently, its subscript is omitted when its dimension is obvious from the context.

Let T_n denote the left-hand matrix in (3.13). \hat{T}_4 is invertible because the matrix A of coefficients of the method is invertible. The inverse of T_n is given by

$$T_n^{-1} = \begin{pmatrix} \hat{T}_1^{-1} + O(h) & O(h) \\ O(h) & \hat{T}_4^{-1}/h + O(1) \end{pmatrix}.$$

Now we can solve for $\tilde{E}^{z'}$ and $\tilde{\tilde{E}}^{z'}$ in (3.13) to obtain

$$\begin{pmatrix}
\tilde{E}^{x'} \\
\tilde{E}^{x'}
\end{pmatrix} = -\begin{pmatrix}
\hat{T}_{1}^{-1}\hat{S}_{1} + O(h) & O(h) \\
O(h) & \hat{T}_{4}^{-1}\hat{S}_{4}/h + O(1)
\end{pmatrix} \begin{pmatrix}
e_{n-1}^{x} + \delta^{x(n)} \\
e_{n-1}^{x} + \delta^{x(n)}
\end{pmatrix} + \begin{pmatrix}
\tilde{\eta}^{x} + O(h)\tilde{\eta}^{y} + O(h)\tilde{\eta}^{x} \\
\hat{T}_{4}^{-1}\tilde{\eta}^{y}/h + O(h)\tilde{\eta}^{x} + O(1)\tilde{\eta}^{y}
\end{pmatrix} (3.14)$$

Recall that $E^{z'} = \tilde{E}^{z'} + \tilde{\tilde{E}}^{z'}$, so that from (3.14),

$$E^{x'} = -(1/h)(\hat{T}_4^{-1}\hat{S}_4)(e_{n-1}^x + \delta^{x(n)}) + (1/h)\hat{T}_4^{-1}\bar{\eta}^y + \bar{\bar{\eta}}^x + O(1)(e_{n-1}^x + \delta^{x(n)}) + O(1)\bar{\eta}^y + O(h)\bar{\bar{\eta}}^x.$$
(3.15)

We can use the expression for $E^{x'}$ above to solve for e_n^x . From the third equation in (3.9), we have

$$e_n^z = e_{n-1}^z + h\mathbf{b}^T E^{z'} + \delta_{M+1}^{z(n)},$$

where $\mathbf{b}^T = (b_1 I_d, b_2 I_d, \dots, b_M I_d) = b^T \otimes I_d$. Substituting (3.15) into the above expression, we obtain

$$e_n^x = e_{n-1}^x - (b^T \otimes I_d)(A^{-1} \otimes I_d)(I_M \otimes H)(e_{n-1}^x + \delta^{z(n)}) + \delta_{M+1}^{z(n)} + (b^T \otimes I_d)(A^{-1} \otimes I_d)\bar{\eta}^y + O(h\delta^{z(n)}) + O(he_{n-1}^x) + O(h\bar{\eta}^y) + O(h\bar{\eta}^z).$$
(3.16)

Now by the strict stability condition, we have $|1 - b^T A^{-1} \epsilon_M| < 1$, where $\epsilon_M = (1, 1, ..., 1)^T$. Let $b^T A^{-1} \epsilon_M = 1 - \rho$. Then from (3.16) we have

$$e_n^x = (I - (1 - \rho)H_{n-1} + O(h))e_{n-1}^x + ((b^T A^{-1}) \otimes I_d)\bar{\eta}^y$$

$$- H_{n-1}((b^T A^{-1}) \otimes I_d)\delta^{x(n)} + \delta_{M+1}^{x(n)}$$

$$+ O(h\delta^{x(n)}) + O(h\bar{\eta}^x) + O(h\bar{\eta}^y). \tag{3.17}$$

Multiplying (3.17) by H_n and using the fact that H_n is a projection, we obtain

$$\tilde{\tilde{e}}_{n}^{x} = \rho(I + O(h))\tilde{\tilde{e}}_{n-1}^{x} - H_{n-1}\left(((b^{T}A^{-1}) \otimes I_{d})\delta^{x(n)} - \delta_{M+1}^{x(n)}\right) + O(h\delta^{x(n)}) + O(h\delta^{x(n)}_{M+1}) + O(\bar{\eta}^{y}) + O(h\bar{\eta}^{x})$$
(3.18)

where $\bar{\eta}^x = (\tilde{\eta}_1^x, \tilde{\eta}_2^x, \dots, \tilde{\eta}_M^x)^T$. Note that, by definition of the algebraic order $k_{a,1}$, we have $((b^T A^{-1}) \otimes I_d) \delta^{x(n)} - \delta^{x(n)}_{M+1} = O(h^{k_{a,1}+1})$. Multiplying (3.17) by $(I - H_n)$, and noting that $(I - H_i) \tilde{\eta}_i^y = 0$ for i = 1, 2, ..., M by definition of $\tilde{\eta}_i^y$ in (3.9), we obtain

$$\tilde{\epsilon}_n^z = (I + O(h))\tilde{\epsilon}_{n-1}^z + O(h\delta^{z(n)}) + O(\delta_{M+1}^{z(n)}) + O(h\bar{\eta}^y) + O(h\bar{\eta}^z). \quad (3.19)$$

Note that the order k_{M+1} of the last stage is always at least as large as the differential order k_d . Now suppose that $\|\bar{\eta}^x\| \leq \epsilon_1$ and $\|\bar{\eta}^y\| \leq \epsilon_2$. The magnitude of ϵ_1 and ϵ_2 will be determined later. For linear problems, they are just proportional to the size of the residuals at the termination of the Newton iteration. Then rewriting (3.18) and (3.19), we have

$$\tilde{\tilde{\epsilon}}_{n}^{x} = \rho(I + O(h))\tilde{\tilde{\epsilon}}_{n-1}^{x} + O(h^{k_{a,1}+1}) + O(h^{k_{I}+2}) + O(\epsilon_{2}) + O(h\epsilon_{1})
\tilde{\epsilon}_{n}^{x} = (I + O(h))\tilde{\epsilon}_{n-1}^{x} + O(h^{k_{I}+2}) + O(h^{k_{d}+1}) + O(h\epsilon_{2}) + O(h\epsilon_{1}), (3.20)$$

where by the strict stability condition, $-1 < \rho < 1$.

Solving the recurrence relations (3.20) and noting that $e_n^z = \tilde{\tilde{e}}_n^z + \tilde{e}_n^z$ and that $k_{a,1} \ge k_I$, we obtain

$$||e_n^z|| = O(h^{k_0}) + O(\epsilon_2) + O(\epsilon_1) + O((I-H)e_0^z),$$
 (3.21)

where $k_G = \min(k_d, k_I + 1)$.

Now we can bound the error in the y component. By the definition of E_i^y ,

$$E^{y} = \mathbf{e}_{n-1}^{y} + h(A \otimes I_{d})E^{y'} + \delta^{y(n)},$$

where $E^{y'} = (E_1^{y'}, E_2^{y'}, \dots, E_M^{y'})^T$, $E^y = (E_1^y, E_2^y, \dots, E_M^y)^T$ and $\delta^{y(n)} = (\delta_1^{y(n)}, \delta_2^{y(n)}, \dots, \delta_M^{y(n)})^T$. We can solve for $E^{y'}$ to obtain

$$E^{y'} = (1/h)(A^{-1} \otimes I_d)(E^y - e_{n-1}^y - \delta^{y(n)}).$$

In the expression for e_n^y given in (3.9), substitute for $E^{y'}$ to obtain

$$e_n^y = (1 - b^T A^{-1} \epsilon_M) e_{n-1}^y - (((b^T A^{-1}) \otimes I_d) \delta^{y(n)} - \delta_{M+1}^{y(n)}) + b^T (A^{-1} \otimes I_d) E^y.$$

Note that $((b^T A^{-1}) \otimes I_d) \delta^{y(n)} - \delta^{y(n)}_{M+1} = O(h^{k_{a,1}+1})$. Substitute for E^y from (3.8) and simplify to obtain

$$e_n^y = \rho e_{n-1}^y - \mathbf{b}^T (A^{-1} \otimes I_d) (I_M \otimes M + O(h)) (E^{x'} + (I_M \otimes G_{11}) E^x) + O(\eta^x) + O(h^{k_{a,1}+1}), \tag{3.22}$$

where M and G_{11} are evaluated along the true solution at t_{n-1} . Now substitute for $E^{z'}$ from (3.15) and simplify, to obtain

$$e_n^y = \rho e_{n-1}^y + (1/h)((b^T(A^{-1})^2) \otimes I_d)(I_M \otimes M)(I_M \otimes H)(e_{n-1}^x + \delta^{x(n)}) + O(e_{n-1}^x) + O(\delta^{x(n)}) + O(\bar{\eta}^y/h) + O(\bar{\eta}^x) + O(h^{k_{a,1}+1}).$$

Noting that $(I_M \otimes H)e_{n-1}^x = (\tilde{\tilde{e}}_{n-1}^x, \tilde{\tilde{e}}_{n-1}^x, \dots, \tilde{\tilde{e}}_{n-1}^x)^T$, and solving the recurrence in (3.20) for $\tilde{\tilde{e}}_{n-1}^x$ and simplifying, we have

$$e_n^y = \rho e_{n-1}^y + O(h^{k_I}) + O(\epsilon_2/h) + O(\epsilon_1).$$

Solving the above recurrence for e_n^y , we obtain

$$||e_n^y|| = O(h^{k_I}) + O(\epsilon_2/h) + O(\epsilon_1). \tag{3.23}$$

Now for a linear problem, if we assume that the residuals from the Newton iteration satisfy $\|\eta_i^x\| = O(h^{k_G})$ and $\|\eta_i^y\| = O(h^{k_I+1})$, then we have shown the desired result for (3.3). For the nonlinear analysis, we will follow

a strategy similar to that used in [14]. Recall that, by definition, η_i^z consists of terms of the form

$$\frac{\partial^2 g_1}{\partial x^2} E_i^x E_i^x, \quad \frac{\partial^2 g_1}{\partial x \partial y} E_i^x E_i^y, \quad \frac{\partial^2 g_1}{\partial y^2} E_i^y E_i^y.$$

while η_i^y is composed of terms of the form

$$\frac{\partial^2 g_2}{\partial x^2} E_i^x E_i^x.$$

Using the solution $\tilde{\tilde{e}}_{n-1}^x$ to (3.20) and equations (3.7), (3.8), (3.15), and (3.21), one can show that $||E_i^y||$ and $||E_i^x||$ satisfy

$$||E_i^y|| = O(h^{k_I}) + O(\epsilon_2/h) + O(\epsilon_1) + O((I-H)\epsilon_0^x) ||E_i^x|| = O(h^{k_O}) + O(\epsilon_2) + O(\epsilon_1) + O((I-H)\epsilon_0^x).$$
(3.24)

If we assume that the residuals from the Newton iteration satisfy $\|\hat{\eta}^x\| = O(h^{k_G})$ and $\|\hat{\eta}^y\| = O(h^{k_I+1})$, where $\hat{\eta}$ is the contribution to η from errors in the Newton iteration, and that $(I-H)e_0^x = O(h^{k_G})$, then we have

$$\begin{aligned} \|\eta^{z}\| & \leq O(h^{k_{G}}) + (O(h^{k_{I}}) + O(\epsilon_{1}) + O(\epsilon_{2}/h))^{2} \\ & \leq O(h^{k_{G}}) + O(h^{2k_{I}}) + O(\epsilon_{1}^{2}) + O(\epsilon_{1}\epsilon_{2}/h) + O(h^{k_{I}}\epsilon_{1}) \\ & + O(h^{k_{I}-1}\epsilon_{2}) + O(\epsilon_{2}^{2}/h^{2}) \\ & \leq K_{1}(h^{k_{G}} + h^{k_{I}}\epsilon_{1} + \epsilon_{1}^{2} + \epsilon_{1}\epsilon_{2}/h + h^{k_{I}-1}\epsilon_{2} + \epsilon_{2}^{2}/h^{2}). \end{aligned}$$

Similarly,

$$\|\eta^{y}\| \leq O(h^{k_{I}+1}) + (O(h^{k_{G}}) + O(\epsilon_{1}) + O(\epsilon_{2}))^{2}$$

$$\leq K_{2}(h^{k_{I}+1} + h^{k_{G}}\epsilon_{1} + \epsilon_{1}^{2} + \epsilon_{1}\epsilon_{2} + h^{k_{G}}\epsilon_{2} + \epsilon_{2}^{2}).$$

Now we let ϵ_1 and ϵ_2 be the solutions to

$$\epsilon_1 = K_1(h^{k_G} + h^{k_I}\epsilon_1 + \epsilon_1^2 + \epsilon_1\epsilon_2/h + h^{k_I-1}\epsilon_2 + \epsilon_2^2/h^2)
\epsilon_2 = K_2(h^{k_I+1} + h^{k_G}\epsilon_1 + \epsilon_1^2 + \epsilon_1\epsilon_2 + h^{k_G}\epsilon_2 + \epsilon_2^2).$$
(3.25)

We would like to conclude that $\epsilon_1 = O(h^{k_G})$, $\epsilon_2 = O(h^{k_I+1})$. Now if we assume that $k_I \geq 2$ and solve (3.25) for ϵ_1 and ϵ_2 by functional iteration, starting with initial values that satisfy $\epsilon_2^{(0)} = O(h^{k_I+1})$ and $\epsilon_1^{(0)} = O(h^{k_G})$, then it is easy to see that the spectral radius of the iteration matrix is less than one, and we can use the contraction mapping theorem to conclude that

it converges to a solution which satisfies $\epsilon_1 = O(h^{k_G})$, $\epsilon_2 = O(h^{k_I+1})$. For $k_I = 1$, we cannot apply the theorem directly because the spectral radius is larger than one, but if we scale the variables by $\bar{\epsilon}_1 = \epsilon_1/\sqrt{h}$ and $\bar{\epsilon}_2 = \epsilon_2/h$, we can then apply the same strategy to reach the conclusion.

We have shown the result for (3.3), and it remains for us to demonstrate that we can extend the conclusions to (3.1). This is easy to do, following arguments similar to those used in [12]. For systems with index one constraints mixed with index two constraints,

$$x' + f_1(x, y, z, t) = 0$$

 $f_2(x, y, t) = 0$
 $f_3(x, t) = 0,$ (3.26)

where $\partial f_2/\partial y$ and $[(\partial f_3/\partial x)(\partial f_1/\partial z)]$ are both nonsingular, we can solve the second equation in (3.26) for Y_i at each stage i, so that the results just shown for the order of x and z are valid. By solving for the error in y_n in terms of the error in y_{n-1} and the internal stage errors in x, it is easy to see that for strictly stable implicit Runge-Kutta methods and consistent initial conditions, the error in y_n is no worse than $O(h^{k_G})$.

Similarly, we can see that the result extends to systems of the form

$$F(x', x, y, z, t) = 0$$

 $f_2(x, y, t) = 0$
 $f_3(x, t) = 0,$ (3.27)

where $\partial F/\partial x'$ is nonsingular, by noting that the first equation in (3.27) can be solved for x' to obtain a system of the form (3.3).

Finally, if in (3.1) $\partial g/\partial y$ is not identically zero but is singular and has constant rank, then we can use a result of Dolezal [8] that there exists smooth nonsingular transformations which bring the system to the form (3.27), and which do not include any change of variables involving x. Thus the conclusions are valid for (3.1).

We should note that this theorem gives only a lower bound on the order of the method, and therefore does not exclude the possibility of a more accurate solution. However, numerical experiments in the next section demonstrate that some implicit Runge-Kutta methods do indeed suffer this order reduction.

4 Numerical Experiments

In this section we present the results of some numerical experiments on linear and nonlinear index two semi-explicit systems. The experiments confirm that the order reduction effects predicted in section 3 can occur in practice, and also raise some interesting questions for future research.

The numerical experiments described in this section were restricted to the four L-stable formulae discussed in section 2. The results given here were obtained using a fixed stepsize code which implements a general M-stage implicit Runge-Kutta method, given the method coefficients. The nonlinear equations at each time step were solved by Newton iteration. The iteration was terminated when the ℓ_2 norm of the difference between two successive iterates was less than a specified tolerance. An analytic iteration matrix was provided to the code for all of the problems. All of the computations were performed in single precision on a CDC 176 computer.

The first test problem was a linear problem having four differential equations and one algebraic equation [2]:

$$x'_{1} = -e^{t}x_{1} + x_{2} + x_{4} + y - e^{-t}$$

$$x'_{2} = -x_{1} + x_{2} - \sin(t)x_{3} + y - \cos(t)$$

$$x'_{3} = \sin(t)x_{1} + x_{3} + \sin(t)x_{4} - \sin^{2}(t) - e^{-t}\sin(t)$$

$$x'_{4} = \cos(t)x_{2} + x_{3} + \sin(t)x_{4} - e^{-t}(1 + \sin(t)) - \cos^{2}(t) - e^{t}$$

$$0 = x_{1}\sin^{2}(t) + x_{2}\cos^{2}(t) + (x_{3} - e^{t})(\sin(t) + 2\cos(t)) + \sin(t)(x_{4} - e^{-t})(\sin(t) + \cos(t) - 1) - \sin^{3}(t) - \cos^{3}(t)$$
(4.1)

The exact solution to this system is $x_1 = \sin(t)$, $x_2 = \cos(t)$, $x_3 = e^t$, $x_4 = e^{-t}$, and $u(t) = e^t \sin(t)$. It is easy to verify that system (4.1) is index two for all t. We solved this test problem for a sequence of fixed stepsizes on the interval [0,1] using the four IRK methods. Consistent initial values were specified at t = 0. After computing the global error at t = 1, an observed rate of convergence was determined by computing the ratio of global errors when successively halving the stepsize. The observed order of the global error was two in all variables when the test problem was solved by the two stage SIRK method, agreeing with the order predicted by the theory. However, when the five stage DIRK was used, we found that the state variables x were computed to an accuracy of $O(h^4)$ (i.e., the nonstiff ODE order k_d), thereby exceeding the lower bound k_G^x on the order predicted in section 3. The algebraic variable y was computed to only O(h) accuracy, which

agrees with the lower bound value of k_G^y . The three stage SIRK method, as expected, determined the algebraic variable to $O(h^2)$ accuracy and the state variables to $O(h^3)$ accuracy. Note that the SIRK methods, as well as the DIRK method, achieved the nonstiff ODE order of accuracy in the state variables for this linear test problem. Finally, we found that the seven stage extrapolation formula was order three in all variables, thereby exceeding the order predicted by the lower bound. From these results, one might be tempted to conclude that the convergence theorem could be strengthened to predict that IRK methods will compute the state variables x to $O(h^{k_d})$ accuracy. However, this is not the case, as we can see from the next two examples. The orders k_g^x and k_g^y observed for this linear test problem are summarized with the predicted lower bounds k_G^x and k_G^y (recall that $k_G^x = \min(k_d, k_I + 1)$ and $k_G^y = k_I$) and the nonstiff ODE and internal orders k_d and k_I in Table 4.1.

Next we investigated the behavior of the IRK formulae on two nonlinear problems. We chose to study the index three pendulum problem simply because it has been studied so frequently by DAE researchers [11],[12],[13] and can be posed as an index two problem [12]. The other nonlinear problem considered arises in the context of trajectory prescribed path control problems [3]. The exact solution is not available for either problem, so we first had to generate a 'true' solution which could be used for comparisons. The corresponding index one systems were formulated and solved by the code DASSL [16] with extremely tight error tolerances. In particular, the 'true' solution to the pendulum problem was obtained by setting the error tolerances RTOL = ATOL = 1.E-12, while it was obtained for the trajectory problem with RTOL = ATOL = 1.E-10.

Consider the pendulum problem as formulated in [12]. Note that this formulation ensures that the original index three algebraic constraint is satisfied even though the index of the system has been reduced to two.

$$x'_{1} = x_{3} - x_{1}y_{2}$$

$$x'_{2} = x_{4} - x_{2}y_{2}$$

$$x'_{3} = -y_{1}x_{1}$$

$$x'_{4} = -y_{1}x_{2} - 1$$

$$0 = (1 - x_{1}^{2} - x_{2}^{2})/2$$

$$0 = x_{1}x_{3} + x_{2}x_{4}$$

$$(4.2)$$

The algebraic constraints in this problem are nonlinear, yet for a constant state the algebraic variables appear only linearly in the system. The pen-

dulum problem was solved using the fixed stepsize IRK code on the interval [0,1] for a sequence of stepsizes with each particular IRK formula. Consistent initial conditions were specified, namely $x_1 = 1$, $x_2 = x_3 = x_4 = y_1 = y_2 = 0$. The corrector iteration was terminated with a tolerance of 10⁻⁸, because the Newton iteration failed to converge for tighter tolerances. Rates of convergence for each method were estimated as in the linear problem by comparing the global errors at t = 1 for numerical solutions produced by successively halving the stepsize. Unlike the results for the linear problem, it does not appear that these formulae determine the state variables x to the ODE order of accuracy. In particular, the five stage DIRK method behaved as expected from the index two convergence theorem, finding the state variables x to no more than $O(h^2)$ accuracy, and the algebraic variables y to O(h) accuracy. Meanwhile, the three stage SIRK method still appeared to be third order in the state variables, while the algebraic variables were determined with close to second order accuracy. The seven stage extrapolation method continued to perform admirably, yielding third order accuracy in all variables. Finally, the two stage SIRK method remained second order accurate for all variables. The numerical results for the pendulum problem are summarized in Table 4.2.

The trajectory problem was posed in [2] as representative of the type of trajectory prescribed path control problems of current interest. The constraints are quite nonlinear in both the state and algebraic variables, while the two algebraic constraints are designed to simply prescribe two of the state variables as functions of time. Initial values for the state variables are known exactly, but initial values for the two algebraic variables (namely, angle of attack α and bank angle β) were determined numerically from the corresponding index one system. Specifically, the test problem used the following initial values for the state variables: altitude H = 100,000 feet, longitude $\xi = 0^{\circ}$, latitude $\lambda = 0^{\circ}$, relative velocity $V_R = 12000$ feet/second, flight path angle $\gamma = -1^{\circ}$, and azimuth $A = 45^{\circ}$. Angle of attack and bank angle were initialized to $\alpha = 2.672870042^{\circ}$ and $\beta = -.0522095861634^{\circ}$, respectively. The 'small' errors in the initial values for the algebraic variables are annihilated in one step by the IRK methods chosen, as a result of their L-stability property. The Newton iteration was terminated with a tolerance of 10^{-10} . This problem was solved for fixed stepsizes on the interval [0,300], and the global errors in the solution were computed at t = 300 using the 'true' solution described earlier. The system is composed of the following six equations of motion and two prescribed path control constraints:

```
H' = V_R \sin(\gamma)
\xi' = V_R \cos(\gamma) \sin(A)/(r \cos(\lambda))
\lambda' = V_R \cos(\gamma) \cos(A)/r
V_R' = -D/m - g \sin(\gamma)
-\Omega_E^2 r \cos(\lambda)(\sin(\lambda) \cos(A) \cos(\gamma) - \cos(\lambda) \sin(\gamma))
\gamma' = L \cos(\beta)/(mV_R) + \cos(\gamma)(V_R^2/r - g)/V_R + 2\Omega_E \cos(\lambda) \sin(A)
+ \Omega_E^2 r \cos(\lambda)(\sin(\lambda) \cos(A) \sin(\gamma) + \cos(\lambda) \cos(\gamma))/V_R
A' = L \sin(\beta)/(mV_R \cos(\gamma)) + V_R \cos(\gamma) \sin(A) \tan(\lambda)/r
-2\Omega_E(\cos(\lambda) \cos(A) \tan(\gamma) - \sin(\lambda))
+ \Omega_E^2 r \cos(\lambda) \sin(\lambda) \sin(A)/(V_R \cos(\gamma))
0 = \gamma + 1 + 9(t/300)^2
0 = A - 45 - 90(t/300)^2
```

where

 $r = H + a_e$

 $a_e = 20902900$ feet, the earth radius $g = \mu/r^2$, the gravity force $\mu = .14076539E + 17$ $\Omega_E = .72921159E - 4$ m = 2.890532728, the mass of the vehicle $L = .5\rho C_L S V_R^2$, the aerodynamic lift force $\rho = .002378e^{-H/23800}$, the atmospheric density $C_L = .01\alpha$, the aerodynamic lift coefficient S = 1, the vehicle cross-sectional reference area $D = .5\rho C_D S V_R^2$, the aerodynamic drag force $C_D = .04 + .1C_L^2$, the aerodynamic drag coefficient.

Since the corrector iteration was terminated with a fairly tight tolerance, the values of the two state variables prescribed by the algebraic constraints (namely, the flight path angle γ and the azimuth A) were computed almost

exactly for all the IRK methods considered. Rates of convergence for all the other variables have been estimated as we described earlier for the other test problems. The results were similar to those obtained for the pendulum problem. The numerical solution produced by the DIRK method was close to second order accurate in the state variables, and first order accurate in the algebraic variables. The extrapolation formula yielded close to third order accurate solutions in all variables, while the two stage SIRK method was clearly second order for all variables. The three stage SIRK method surprised us by producing a third order accurate solution for the algebraic variables as well as for the state variables. We suspect this difference in performance for this particular IRK method, when compared to its results on the linear test problem and the pendulum problem, must be due to the specific coupling of the state and algebraic variables in this nonlinear system. The numerical results for the trajectory problem are summarized in Table 4.3.

In conclusion, we see that the observed convergence rates of these IRK methods applied to nonlinear semi-explicit index two systems can sometimes be as slow as the lower bounds derived in section 3 would indicate. Some formulae, in particular the extrapolation method, achieve an order of accuracy exceeding the predicted lower bounds, suggesting that a stronger convergence theorem might be possible. On the other hand, there is a class of IRK methods which have an internal order as high or nearly as high as the ODE order. In particular, consider the class of M-stage singly-implicit Runge-Kutta methods (SIRKs) whose coefficient matrix A is characterized by its single-fold eigenvalue. Butcher [6] has shown how these IRK formulae can be implemented very efficiently. There are two types of SIRKs, the transformed type [5] and the collocation type [15]. It is easy to show that for index two problems, transformed SIRKs will be at least order M-1 (since $k_I \geq M-1$), while collocation SIRKs will be order M (since $k_I = M$). Note that the second order, two-stage SIRK formula which has appeared so promising in our numerical experiments is in fact a collocation method. Note also that if an L-stable SIRK formula is desired, one may select the eigenvalue of the A matrix to satisfy $L_M(\lambda^{-1}) = 0$ where L_M is the Laquerre polynomial of degree M. Methods of this type have been derived for orders up to and including six. In summary, we expect the SIRK methods to perform very well on index two problems. However, the development of an efficient IRK code for DAEs with index greater than one remains a challenge because of the difficulties in developing appropriate error control strategies for all the variables.

Table 4.1: Predicted/Observed Orders for Linear Test Problem

L-stable Methods	k_d	k_I	k_G^x	k_g^x	k_G^y	k_g^y
1. Two-stage SIRK	2	2	2	2	2	2
2. Five-stage DIRK	4	1	2	4	1	1
3. Three-stage SIRK	3	2	3	3	2	2
4. Seven-stage Extrp.	3	1	2	3	1	3

Table 4.2: Predicted/Observed Orders on Pendulum Problem

L-stable Methods		k_I	k_G^x	k_q^x	k_G^y	k_a^y
1. Two-stage SIRK	2	2	2	2	2	2
2. Five-stage DIRK	4	1	2	2	1	1
3. Three-stage SIRK	3	2	3	3	2	2
4. Seven-stage Extrp.	3	1	2	3	1	3

Table 4.3: Predicted/Observed Orders on Trajectory Problem

L-stable Methods		k _I	k_G^z	k_g^x	k_G^y	k_g^y
1. Two-stage SIRK	2	2	2	2	2	2
2. Five-stage DIRK	4	1	2	2	1	1
3. Three-stage SIRK	3	2	3	3	2	3
4. Seven-stage Extrp.	3	1	2	3	1	3

References

- [1] R. Alexander, Diagonally Implicit Runge-Kutta Methods for Stiff ODEs, SIAM J. Numer. Anal. 14 (1977), pp. 1006-1022.
- [2] K. E. Brenan, Stability and Convergence of Difference Approximations for Higher Index Differential-Algebraic Systems with Applications in Trajectory Control, Ph.D. Thesis, University of California at Los Angeles, 1983.
- [3] K. E. Brenan, Numerical Simulation of Trajectory Prescribed Path Control Problems by the Backward Differentiation Formulas, IEEE Journal on Automatic Control, AC-31 (1986), pp. 266-269.
- [4] K. Burrage, Efficiently Implementable Algebraically Stable Runge-Kutta Methods, SIAM J. Numer. Anal. 19 (1982), pp. 245-258.
- [5] K. Burrage, A Special Family of Runge-Kutta Methods for Solving Stiff Differential Equations, BIT 18 (1978), pp. 22-41.
- [6] J. C. Butcher, On the Implementation of Runge-Kutta Methods, BIT 16 (1976), pp. 237-240.
- [7] J. R. Cash, Diagonally Implicit Runge-Kutta Formulae with Error Estimates, J. Inst. Maths. Applics., 24 (1979), pp. 293-301.
- [8] V. Dolezal, The Existence of a Continuous Basis of a Certain Linear Subspace of E_r, which Depends on a Parameter, Časopis pro pěstováni matematiky, roč. 89, Praha, (1964), pp. 466-468.
- [9] R. Frank, J. Schneid and C. W. Ueberhuber, Order Results for Implicit Runge-Kutta Methods Applied to Stiff Systems, SIAM J. Numer. Anal. 22 (1985), pp. 515-534.
- [10] F. R. Gantmacher, The Theory of Matrices, Vol. 2, Chelsea, New York, 1964.
- [11] C. W. Gear, H. H. Hsu, and L. Petzold, Differential-Algebraic Equations Revisited, Proc. Conference on Matrix Pencils, Piteå, Sweden, 1982.
- [12] C. W. Gear, B. Leimkuhler and G. K. Gupta, Automatic Integration of Euler-Lagrange Equations with Constraints, Journal of Computational and Applied Mathematics 12 and 13 (1985), pp. 77-90.

- [13] B. J. Leimkuhler, Error Estimates for Differential-Algebraic Equations, Report No. UIUCDCD-R-86-1287, Dept. of Computer Science, Univ. of Illinois, 1986.
- [14] P. Lötstedt and L. Petzold, Numerical Solution of Nonlinear Differential Equations with Algebraic Constraints I: Convergence Results for Backward Differentiation Formulas, Mathematics of Computation 46 (1986), pp. 491-516.
- [15] Nørsett, S.P., Runge-Kutta Methods with a Multiple Real Eigenvalue Only, BIT 16 (1976), pp. 388-393.
- [16] L. R. Petzold, A Description of DASSL: A Differential/Algebraic System Solver, IMACS Trans. on Scientific Computation, Vol. 1, R. S. Stepleman ed., 1982.
- [17] L. R. Petzold, Order Results for Implicit Runge-Kutta Methods Applied to Differential/Algebraic Systems, SIAM J. Numer. Anal. 23 (1986), pp. 837-852.
- [18] A. Prothero and A. Robinson, On the Stability and Accuracy of One-Step Methods for Solving Stiff Systems of Ordinary Differential Equations, Math. Comp. 28 (1974), pp. 145-162.